



Large-Scale **Adaptive Semi-Supervised** Learning via Unified Inductive and **Transductive** Model 1-KDD 2014

Active-Transductive Learning with **Label-Adapted Kernels** 2-KDD 2014

Semisupervised Dimensionality Reduction and Classification Through Virtual Label Regression 3-IEEE Transactions on ...Part B (Cybernetics) 2011

Transductive Learning on **Adaptive Graphs** 4-AAAI 2010

Transductive Classification via Local Learning Regularization 5-AISTAS 2007

Linear **Manifold Regularization** with **Adaptive Graph** for **Semi-supervised** Dimensionality Reduction 6-IJCAI 2017

Semi-supervised manifold regularization with **adaptive graph** construction 7-PRL 2017



Active Learning via Transductive Experimental Design 8-ICML 2006

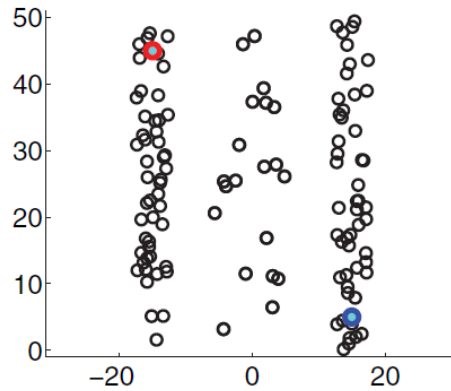
Non-greedy Active Learning for Text Categorization using Convex Transductive Experimental Design 9-SIGIR 2008

Beyond the Point Cloud: from Transductive to Semi-supervised Learning 10-ICML 2005

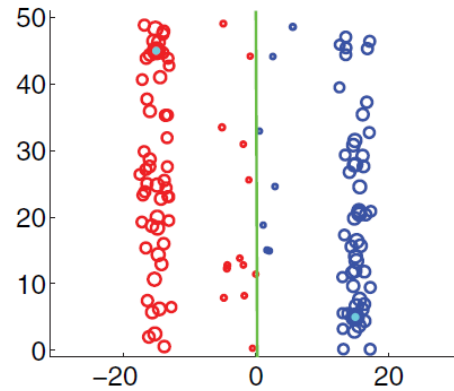
Manifold Regularized Experimental Design for Active Learning 11-TIP 2017

Diversifying Convex Transductive Experimental Design for Active Learning 12-IJCAI 2016

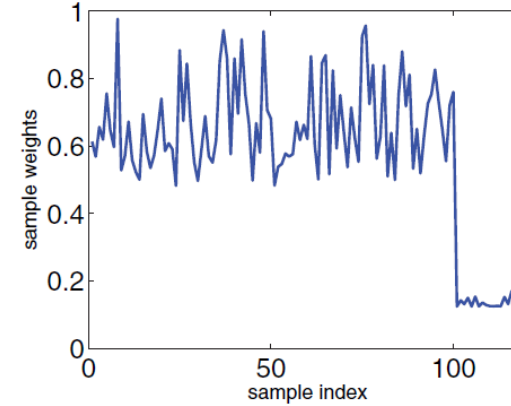
1-KDD 2014



(a) Original toy data



(b) Toy data after classification using our model

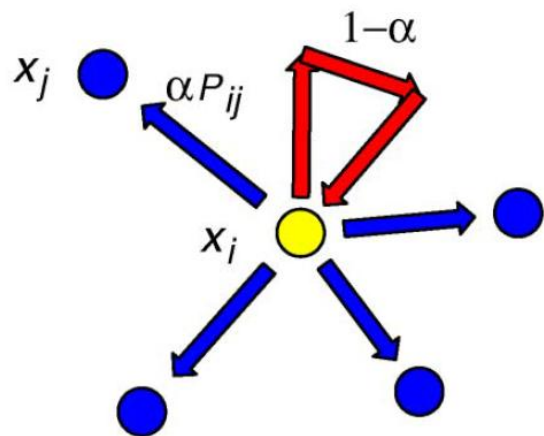


(c) Sample weights

3-TCYB 2011

In order to capture the outliers or the data from the class other than the labeled classes, we consider special random walks

$$\mathbf{A}_{ij} = e^{-\|x_i - x_j\|^2 / \sigma^2} \quad \mathbf{P} = \mathbf{D}^{-1} \mathbf{A} \quad \mathbf{D}_{ii} = \sum_j \mathbf{A}_{ij}$$



$$\tilde{\mathbf{P}} = \mathbf{I}_\beta + \mathbf{I}_\alpha \mathbf{P} \quad \mathbf{I}_\beta = \mathbf{I} - \mathbf{I}_\alpha$$

$$\mathbf{Y} = [\mathbf{Y}_1^T, \dots, \mathbf{Y}_n^T]^T \in \mathbb{R}^{n \times (c+1)}$$

add an additional class $c + 1$ in order to detect outlier data

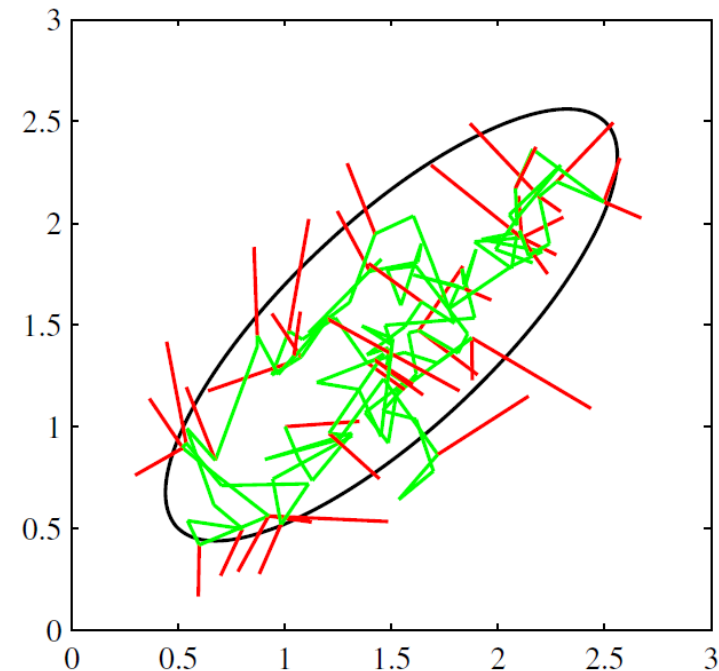
$$\mathbf{G} = \mathbf{I}_\beta + \hat{\mathbf{P}}\mathbf{I}_\beta + \hat{\mathbf{P}}^2\mathbf{I}_\beta + \dots + \hat{\mathbf{P}}^n\mathbf{I}_\beta + \dots \quad \hat{\mathbf{P}} = \mathbf{I}_\alpha \mathbf{P}$$

$$\mathbf{G} = (\mathbf{I} - \mathbf{I}_\alpha \mathbf{P})^{-1} \mathbf{I}_\beta$$

the probability of the i_{th} point stopping the walks at the j_{th} point.

$$\mathbf{F} = \mathbf{G}\mathbf{Y}$$

the probability of the i_{th} point which stops the random walks at the labeled data point whose label is j





In the special random walks, if an unlabeled sample **is similar to one of the labeled samples**, the walks starting from this unlabeled sample will **stop at the labeled sample with high probability**.

If this unlabeled **sample is an outlier or from a novel class**, this sample is not similar to any of the labeled samples, and the walks starting from this unlabeled sample will **stop at one of the unlabeled samples with high probability**.

Label Regression

$$\mathbf{y} = \mathbf{W}^T \mathbf{x} + \mathbf{b} \quad \mathbf{t}_i = [\underbrace{0, \dots, 0}_{i-1}, 1, \underbrace{0, \dots, 0}_{c-i}]^T$$

$$\mathbf{W} = \arg \min \mathcal{J}(\mathbf{W}, \mathbf{b}) \quad \mathcal{J}(\mathbf{W}, \mathbf{b}) = \gamma \|\mathbf{W}\|^2 + \sum_{i=1}^n \sum_{j=1}^c \mathbf{F}_{ij} \|\mathbf{W}^T \mathbf{x}_i + \mathbf{b} - \mathbf{t}_j\|^2$$

Adaptive Semi-Supervised Learning

$$\left(\begin{array}{l} \min_{\mathbf{W}, \mathbf{b}, \mathbf{Y}} \left\| \mathbf{X}_l^T \mathbf{W} + \mathbf{1}_{nl} \mathbf{b}^T - \mathbf{Y}_l \right\|_F^2 + \sum_{i=1}^n \sum_{k=1}^c y_{ik}^r \left\| \mathbf{x}_i^T \mathbf{W} + \mathbf{b}^T - \mathbf{t}_k^T \right\|_F^2 \\ s.t. \quad \forall i, y_{ik} \in [0, 1], \sum_{k=1}^c y_{ik} = 1 \end{array} \right) \longrightarrow \text{1-KDD 2014}$$



For clearly classified points, y_{i1}, y_{i2} would be one large and one small, For boundary points, however, y_{i1}, y_{i2} would be more likely equal

$$y_{i1} = 0.9, y_{i2} = 0.1, y_{i1}^r = 0.81, y_{i2}^r = 0.01 \quad y_{i1} = 0.5, y_{i2} = 0.5, y_{i1}^r = 0.25, y_{i2}^r = 0.25$$

2-KDD 2014

Classification via Label-Adapted Kernels

Definition 1. The label-adapted weight matrix is defined as

$$W^{\bar{x}}(x_i, x_j) = m_1 \left(\frac{\rho_1(x_i, x_j)^2}{\sigma_1} \right) m_2 \left(\frac{\rho_2(\bar{\chi}(x_i), \bar{\chi}(x_j))^2}{\sigma_2} \right),$$

In particular, the minimization of the quadratic energy relies on the assumption that data points that belong to different classes will have a low similarity weight. This smoothness assumption, which is used in semi-supervised classification tasks, is often both local and global.

However, in many real data sets this might not be the situation, as h may not be smooth with respect to W



Transductive Classification (TC) problem: given labeled data $(x_1, y_1), \dots, (x_l, y_l)$ and unlabeled data $x_{l+1}, \dots, x_{l+u}, y_i \in \{-1, +1\}, 1 \leq i \leq l$. The goal is to predict the class labels of the given unlabeled data.

$$\min_{\mathbf{f} \in \mathbb{R}^n} \mathbf{f}^T \mathbf{L} \mathbf{f} + (\mathbf{f} - \mathbf{y})^T \mathbf{C} (\mathbf{f} - \mathbf{y}) \quad \mathbf{f} = (f_1, \dots, f_n)^T \in \mathbb{R}^n \quad \mathbf{L} = \mathbf{D} - \mathbf{W}$$

$\mathbf{C} \in \mathbb{R}^{n \times n}$ is a diagonal matrix with the i_{th} diagonal element c_i
 set as: $c_i = \alpha_l > 0 \quad 1 \leq i \leq l \quad c_i = \alpha_u \geq 0 \quad l+1 \leq i \leq l+u$

$$w_{ij} = \exp \left(- \frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2} \right)$$

The first term of the objective function evaluates **how smooth the prediction \mathbf{f} is with respect to the data manifold** which is represented by the regularization matrix \mathbf{L} . It reflects the prior knowledge that a good prediction should **have low variance along the data manifold**.

The second term of the objective function uses a quadratic loss function to **measure the fitting error** of the prediction \mathbf{f} . It constrains \mathbf{f} to be close to \mathbf{y} .



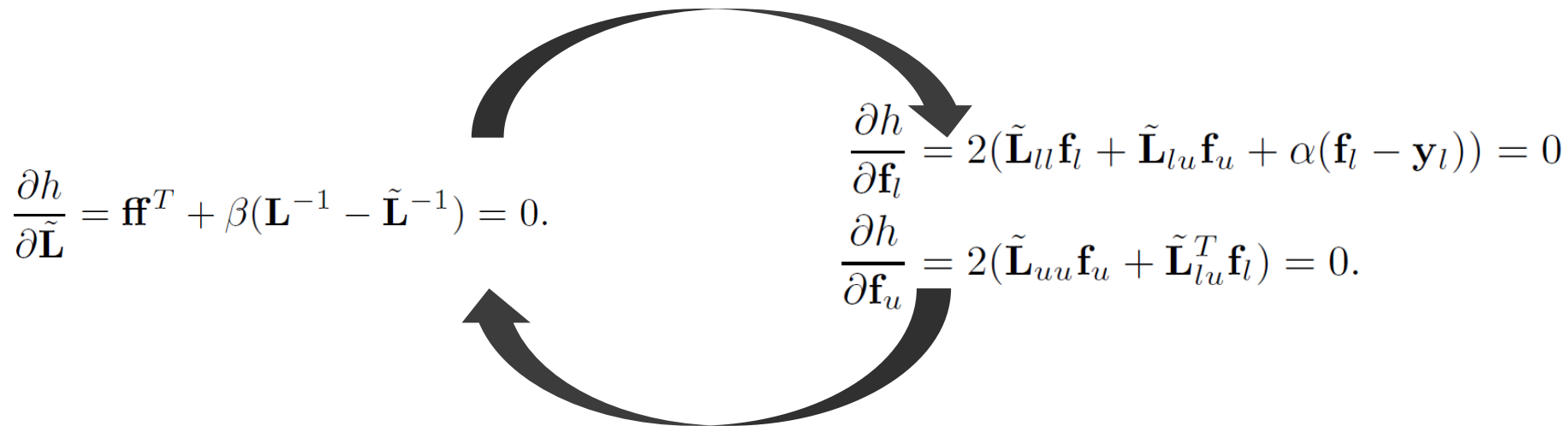
GRF: f_i must be strictly equal to y_i for $1 \leq i \leq l$ and there is no constraint on the unlabeled data.

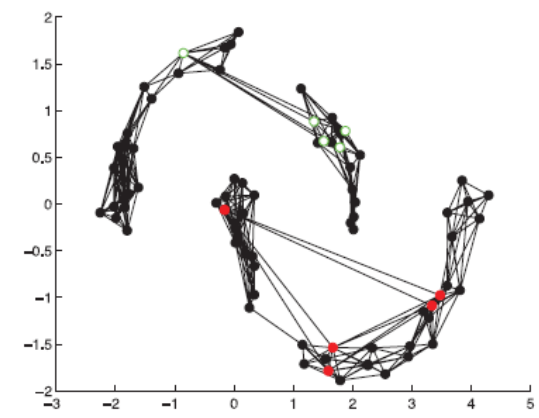
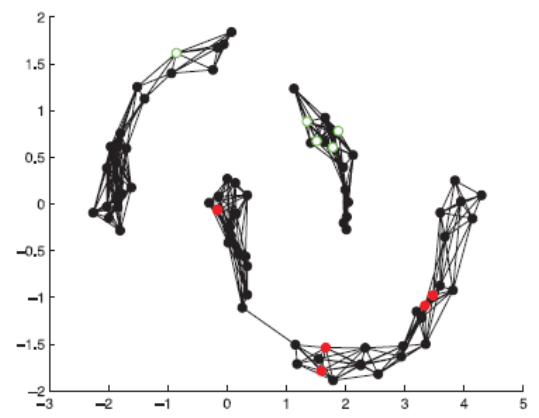
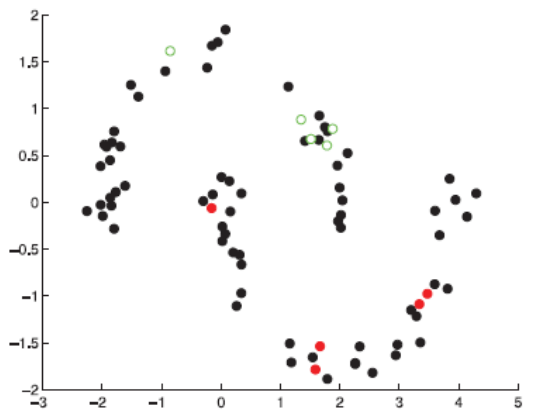
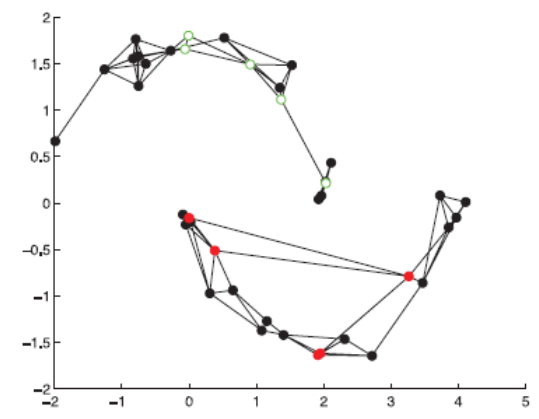
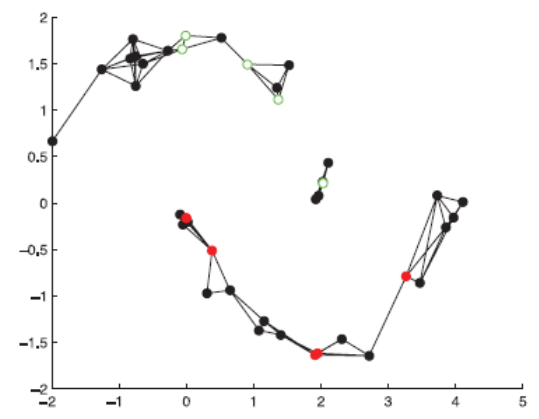
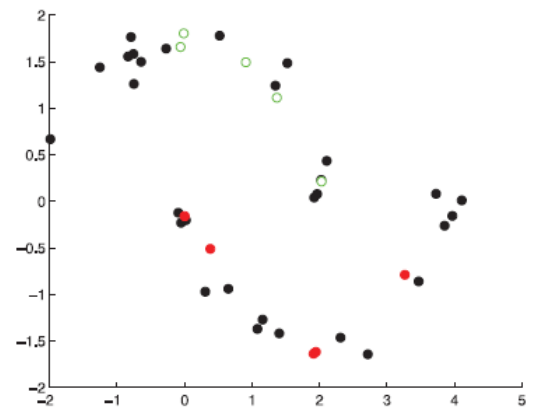
LapRLS: imposes a soft constraint on the labeled data but no constraint on the unlabeled data.

4-AAAI 2010

$$\min_{\mathbf{f} \in \mathbb{R}^n, \tilde{\mathbf{L}} \succeq 0} h(\mathbf{f}, \tilde{\mathbf{L}}) = \mathbf{f}^T \tilde{\mathbf{L}} \mathbf{f} + (\mathbf{f} - \mathbf{y})^T \mathbf{C} (\mathbf{f} - \mathbf{y}) + \beta D(\tilde{\mathbf{L}}, \mathbf{L})$$

$D(A,B)$ is a measure of the dissimilarity between two matrices







JMLR 2006

$$\text{LapRLS} \quad \min_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l (y_i - f(x_i))^2 + \gamma_A \|f\|_K^2 + \frac{\gamma_I}{(u+l)^2} \mathbf{f}^T L \mathbf{f}$$

$$\text{LapSVM} \quad \min_{f \in \mathcal{H}_K} \frac{1}{l} \sum_{i=1}^l (1 - y_i f(x_i))_+ + \gamma_A \|f\|_K^2 + \frac{\gamma_I}{(u+l)^2} \mathbf{f}^T L \mathbf{f}$$

6-IJCAI 2017

Linear **Manifold Regularization** with **Adaptive Graph** for **Semi-supervised** Dimensionality Reduction

$$\min_{W, b} \gamma_A \|W\|_F^2 + \gamma_I \text{Tr}(W^T X L X^T W) + \frac{1}{l} \sum_{i=1}^l \|W^T x_i + b - y_i\|^2$$

Adaptive Neighbor Learning

$$\min_{s_i^T \mathbf{1}=1, 0 \leq s_{ij} \leq 1} \sum_{i,j=1}^n (\|x_i - x_j\|_2^2 s_{ij} + \gamma s_{ij}^2) \quad \text{Tr}(W^T X L X^T W) = \frac{1}{2} \sum_{i,j=1}^n \|W^T x_i - W^T x_j\|_2^2 s_{ij}$$

$$\min_{W, b, S} \sum_{i,j=1}^n \|W^T x_i - W^T x_j\|_2^2 s_{ij} + \gamma \|S\|_F^2 + \beta \|W\|_F^2 + \alpha \text{Tr}(W^T X + b \mathbf{1}^T - Y) U (W^T X + b \mathbf{1}^T - Y)^T$$

$$s.t. \quad S \geq 0, S^T \mathbf{1} = \mathbf{1}$$



To control the uniformity level of the manifold graph weights, we use an entropy regularization term.

$$\min_{f, w_{ij}} \sum_{i,j=1}^{l+u} w_{ij} (f(x_i) - f(x_j))^2 + \gamma_1 \sum_{i=1}^l (f(x_i) - y_i)^2 + \gamma_2 \|f\|_K^2 + \eta \sum_{i,j=1}^n w_{ij} \ln w_{ij}$$

$$\text{s. t. } \sum_{j=1}^u w_{ij} = 1, w_{ij} > 0$$

In fact, each instance in the manifold graph is connected with only a few neighbor instances, thus only a few elements of each weight w_i should be non-zeros, and the rests should be zeros. That is, each vector weight vector w_i **should be sparse**.

$$\min_{f, w_{ij}} \sum_{i,j=1}^{l+u} w_{ij} (f(x_i) - f(x_j))^2 + \gamma_1 \sum_{i=1}^l (f(x_i) - y_i)^2 + \gamma_2 \|f\|_K^2 + \eta_1 \sum_{i=1}^u (x_i - \sum_{i=1}^u w_{ij} x_j)^2 + \eta_2 \sum_{i=1}^{l+u} \|w_i\|_1$$

$$\text{s. t. } \sum_{j=1}^u w_{ij} = 1, w_{ij} > 0$$



Active Learning via Transductive Experimental Design 8-ICML 2006

Non-greedy Active Learning for Text Categorization using Convex Transductive Experimental Design 9-SIGIR 2008

Beyond the Point Cloud: from Transductive to Semi-supervised Learning 10-ICML 2005

Manifold Regularized Experimental Design for Active Learning 11-TIP 2017

Diversifying Convex Transductive Experimental Design for Active Learning 12-IJCAI 2016



Experimental design

Classic experiment design considers learning a linear function $f(x) = w^T x$ $w \in \mathbb{R}^d$, from measurements $y_i = w^T x_i + \epsilon_i, i = 1, \dots, m, \epsilon_i \sim N(0, \sigma^2)$ is measurement noise. x_1, \dots, x_m **are are experiments chosen from n candidates** $v_1, \dots, v_n \in \mathbb{R}^d, n > m$.

The goal of experimental design is to **find a set of experiments x_i that together are maximally informative.**

$$X: [x_1, \dots, x_m]^T \in \mathbb{R}^{m \times d}, \text{set}\{x_i\} \quad |X| = m \quad V: [v_1, \dots, v_n]^T \in \mathbb{R}^{n \times d}, \text{set}\{v_i\} \quad |V| = n$$

The maximum-likelihood estimate of w :

$$\hat{w} = \arg \min_{\mathbf{w}} \left\{ J(\mathbf{w}) = \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 \right\}$$



Estimation error $e = w - \hat{w}$ mean: 0 covariance matrix: $\sigma^2 C_w$, C_w is the inverted Hessian of $J(w)$, σ is a constant

$$C_w = \left(\sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} = (\mathbf{X}^\top \mathbf{X})^{-1}$$

The matrix C_w **characterizes the confidence of the estimation, or the informativeness of the selected data.**

A-optimal design

Let m_j denote the number of times for which v_j is chosen in \mathbf{X} , $m_1 + \dots + m_n = m$

$$\min_{m_1, \dots, m_n} \text{Tr} \left[\left(\sum_{j=1}^n m_j \mathbf{v}_j \mathbf{v}_j^\top \right)^{-1} \right]$$

subject to $m_j \geq 0, m_1 + \dots + m_n = m, m_i \in \mathbb{Z}$

$$\tau_j = m_j / m$$

$$\min_{\tau_1, \dots, \tau_n} \text{Tr} \left[\left(\sum_{j=1}^n \tau_j \mathbf{v}_j \mathbf{v}_j^\top \right)^{-1} \right] \quad \text{subject to } \tau \succeq 0, \mathbf{1}^\top \tau = 1$$

Transductive Experimental Design



- ✦ The optimization criteria based on C_w **does not directly characterize the quality of predictions** on test data
- ✦ Standard experimental design **only considers linear functions** and is thus restrictive in applications.
- ✦ Very importantly, classic experimental design has to solve a SDP problem, which is often **very slow** when dealing with hundreds of data points

A general setting may consider a different set \mathbf{T} of test data points besides candidates in \mathbf{V} . Assume the two sets are the same.

$$\min_{\mathbf{w}} \left\{ J(\mathbf{w}) = \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \mu \|\mathbf{w}\|^2 \right\}$$

$$\mathbf{C}_w = \left(\frac{\partial J(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right)^{-1} = (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1}$$



$\mathbf{f} = [f(v_1), \dots, f(v_n)]$ be the function values on all the available data \mathbf{V} ,
the predictive error $\mathbf{f} - \hat{\mathbf{f}}$ has the covariance matrix $\sigma^2 \mathbf{C}_f$

$$\begin{aligned} \mathbf{C}_f &= \mathbf{V} \mathbf{C}_w \mathbf{V}^\top = \mathbf{V} (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1} \mathbf{V}^\top \\ &= \frac{1}{\mu} \left[\mathbf{V} \mathbf{V}^\top - \mathbf{V} \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \mu \mathbf{I})^{-1} \mathbf{X} \mathbf{V}^\top \right] \end{aligned}$$

Woodbury matrix identity

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{C}^{-1} + \mathbf{V} \mathbf{A}^{-1} \mathbf{U})^{-1} \mathbf{V} \mathbf{A}^{-1}$$

The average predictive variance on \mathbf{V} is given by $\frac{\sigma^2}{n} \text{Tr}(\mathbf{C}_f)$

Transductive experimental design

$$\begin{aligned} &\max_{\mathbf{X}} \text{Tr} \left[\mathbf{V} \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \mu \mathbf{I})^{-1} \mathbf{X} \mathbf{V}^\top \right] \\ &\text{subject to } \mathbf{X} \subset \mathbf{V}, |\mathbf{X}| = m \end{aligned}$$

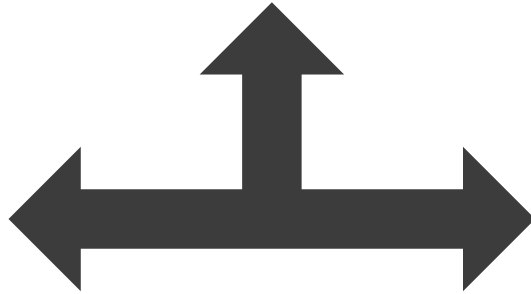
$$\text{Tr}(\mathbf{C}_f) = \text{Tr}(\mathbf{C}_w \mathbf{V}^\top \mathbf{V})$$



Kernel Transductive Experimental Design

$$k(\mathbf{x}, \mathbf{v}) = \langle \phi(\mathbf{x}), \phi(\mathbf{v}) \rangle, \quad \mathbf{x}, \mathbf{v} \in \mathbb{R}^d$$

$$\begin{aligned} & \max_{\mathbf{X}} \quad \text{Tr} [\mathbf{K}_{\mathbf{v}\mathbf{x}}(\mathbf{K}_{\mathbf{x}\mathbf{x}} + \mu\mathbf{I})^{-1}\mathbf{K}_{\mathbf{x}\mathbf{v}}] \\ \text{subject to} \quad & \mathbf{X} \subset \mathbf{V}, |\mathbf{X}| = m \end{aligned}$$



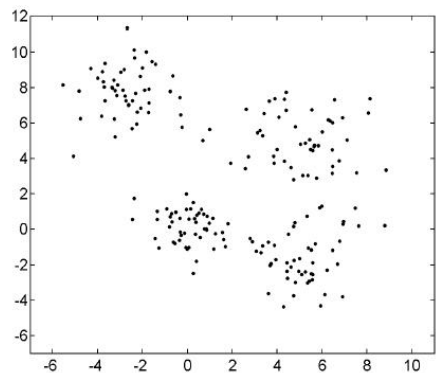
$$\begin{aligned} & \max_{\mathbf{X}} \quad \text{Tr} [\mathbf{V}\mathbf{X}^\top(\mathbf{X}\mathbf{X}^\top + \mu\mathbf{I})^{-1}\mathbf{X}\mathbf{V}^\top] \\ \text{subject to} \quad & \mathbf{X} \subset \mathbf{V}, |\mathbf{X}| = m \end{aligned}$$

$$(\mathbf{K})_{ij} = k(\mathbf{v}_i, \mathbf{v}_j), \quad (\mathbf{K}_{\mathbf{v}\mathbf{x}})_{ij} = k(\mathbf{v}_i, \mathbf{x}_j) \quad (\mathbf{K}_{\mathbf{x}\mathbf{x}})_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

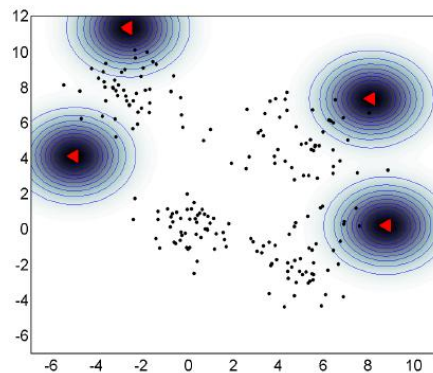
Sequential Design

- Select $\mathbf{x} \in \mathbf{V}$ with the highest $\|\mathbf{K}_{\mathbf{x}}\|^2 / (k(\mathbf{x}, \mathbf{x}) + \mu)$, and add \mathbf{x} into \mathbf{X} , where $\mathbf{K}_{\mathbf{x}}$ and $k(\mathbf{x}, \mathbf{x})$ are \mathbf{x} 's corresponding column and diagonal entry in current \mathbf{K} ;

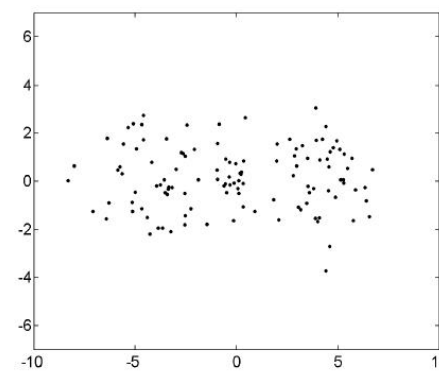
- Update $\mathbf{K} \leftarrow \mathbf{K} - \frac{\mathbf{K}_{\mathbf{x}}\mathbf{K}_{\mathbf{x}}^\top}{(k(\mathbf{x}, \mathbf{x}) + \mu)}$



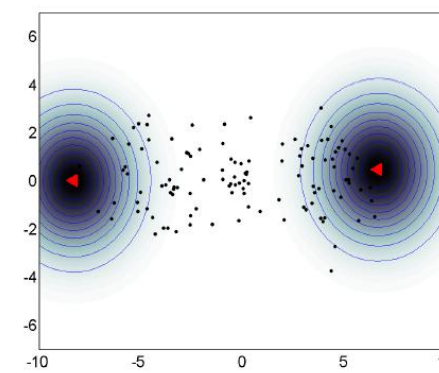
(a) Data set



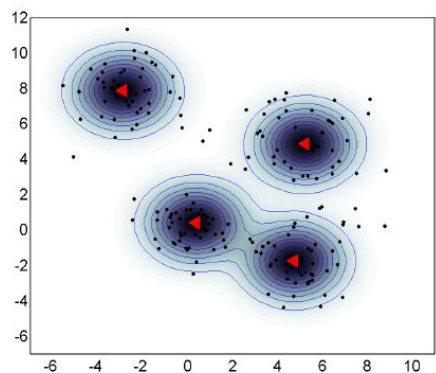
(b) A-optimal design



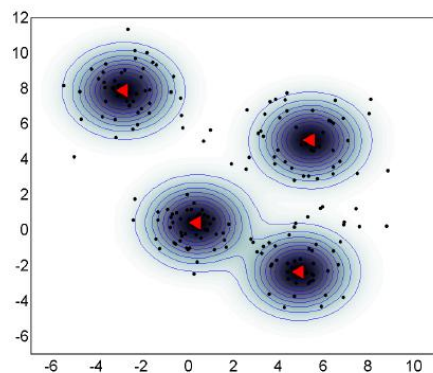
(a) Data set



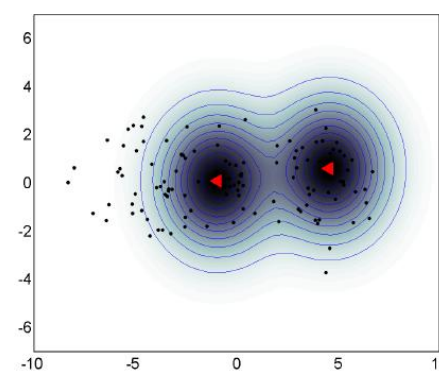
(b) A-optimal design



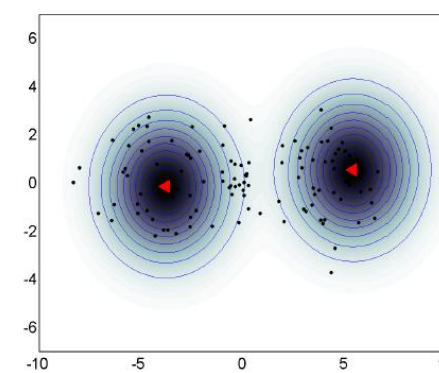
(c) Sequential design



(d) Alternating design



(c) Sequential design



(d) Alternating design



$$\min_{\boldsymbol{\beta}, \boldsymbol{\alpha}_i \in \mathbb{R}^N} \sum_{i=1}^M \|\mathbf{x}_i - \mathbf{X}_C^T \mathbf{B} \boldsymbol{\alpha}_i\|^2 + \mu \|\mathbf{B} \boldsymbol{\alpha}_i\|^2 + \gamma \|\boldsymbol{\beta}\|_1$$

subject to $\mathbf{x}_i \in \mathbf{X}_P, \mathbf{B} = \text{diag}(\boldsymbol{\beta}), \mathbf{B} \succeq 0.$

9-SIGIR 2008

$$\min_{\boldsymbol{\beta}, \boldsymbol{\alpha}_i \in \mathbb{R}^N} \sum_{i=1}^M \|\mathbf{x}_i - \mathbf{X}_C^T \boldsymbol{\alpha}_i\|^2 + \sum_{j=1}^N \frac{\alpha_{i,j}^2}{\beta_j} + \gamma \|\boldsymbol{\beta}\|_1$$

subject to $\mathbf{x}_i \in \mathbf{X}_P, \beta_j \geq 0, j = 1, \dots, N,$

$$\boldsymbol{\beta} = [\beta_1, \dots, \beta_N]$$

$$\boldsymbol{\alpha}_i = [\alpha_{i,1}, \dots, \alpha_{i,N}]^T$$

$$\frac{\sum_{i=1}^M \alpha_{i,j}^2}{|\beta_j|} + \gamma |\beta_j| \geq 2 \sqrt{\gamma \sum_{i=1}^M \alpha_{i,j}^2}, \quad \beta_j^2 = \frac{1}{\gamma} \sum_{i=1}^M \alpha_{i,j}^2$$

$$\min_{\boldsymbol{\alpha}_i \in \mathbb{R}^N} \sum_{i=1}^M \|\mathbf{x}_i - \mathbf{X}_C^T \boldsymbol{\alpha}_i\|^2 + 2 \sum_{j=1}^N \sqrt{\gamma \sum_{i=1}^M \alpha_{i,j}^2}$$

Find the informative samples by **minimizing the expected prediction variance** on the test data.

$$\tilde{k}(x, z) = k(x, z) - \mathbf{k}_x^t (I + MK)^{-1} M \mathbf{k}_z$$

10-ICML 2005

$$\min_{\alpha_i \in \mathbb{R}^l} \sum_{i=1}^l \|\tilde{\phi}(x_i) - \tilde{\phi}(Z) \alpha_i\|^2 + \gamma_1 \|\alpha_i\|^2$$

$$\tilde{\phi}(Z) = [\tilde{\phi}(z_1), \dots, \tilde{\phi}(z_l)]$$

$$\min_{\alpha_i, \beta \in \mathbb{R}^l} \sum_{i=1}^l \left(\|\tilde{\phi}(x_i) - \tilde{\phi}(X) \alpha_i\|^2 + \sum_{j=1}^n \frac{\alpha_{i,j}^2}{\beta_j} \right) + \lambda \|\beta\|_1$$

s. t. $\beta_j \geq 0, \quad j = 1, \dots, n$

$$\alpha_i = (\alpha_{i,1}, \dots, \alpha_{i,n})^T$$

分

割

线

$$\min_{\mathcal{A}, \alpha_i \in \mathbb{R}^K} \sum_{i=1}^M \|\mathbf{x}_i - \mathbf{X}_{\mathcal{A}}^T \alpha_i\|^2 + \mu \|\alpha_i\|^2$$

subject to $|\mathcal{A}| = K, \quad \mathcal{A} \subset \mathcal{C}, \quad \mathbf{x}_i \in \mathbf{X}_{\mathcal{P}}$

$$\min_{\beta, \alpha_i \in \mathbb{R}^N} \sum_{i=1}^M \|\mathbf{x}_i - \mathbf{X}_{\mathcal{C}}^T \alpha_i\|^2 + \sum_{j=1}^N \frac{\alpha_{i,j}^2}{\beta_j} + \gamma \|\beta\|_1$$

subject to $\mathbf{x}_i \in \mathbf{X}_{\mathcal{P}}, \quad \beta_j \geq 0, \quad j = 1, \dots, N,$

$$\beta = [\beta_1, \dots, \beta_N]$$

$$\alpha_i = [\alpha_{i,1}, \dots, \alpha_{i,N}]^T$$



9-SIGIR 2008



11-TIP 2017



Algorithm 2 Convex TED

Require: candidates \mathbf{X}_C , unlabeled data \mathbf{X}_P , $\gamma > 0$;

1: initialize $(\alpha_{i,j})$;

2: **repeat**

3: $\beta_j \leftarrow \sqrt{\frac{1}{\gamma} \sum_{i=1}^M \alpha_{i,j}^2}$ for $j = 1, \dots, N$;

4: $\alpha_i \leftarrow (\text{diag}(\beta)^{-1} + \mathbf{X}_C \mathbf{X}_C^\top)^{-1} \mathbf{X}_C \mathbf{x}_i$, for $i = 1, \dots, M$;

5: **until** converge;

6: $\mathbf{X}_A \leftarrow \{\mathbf{x}_j | \mathbf{x}_j \in \mathbf{X}_C, \beta_j \neq 0\}$;

7: **return** \mathbf{X}_A

Algorithm 1 MRED for Active Learning

Input: The n unlabeled data samples X , the number of the selected most information data samples l , the number of the nearest neighbor data samples k

Step 1: Construct a nearest neighbor Laplacian graph with the weight matrix W as calculated in Eq. (15) on the unlabeled samples X and calculate

Step 2: Construct the kernel Gram matrix K with an selected input kernel type and let $M = L$.

Step 3: Construct the data-dependent deformed kernel Gram matrix \tilde{K} according to Eq. (14).

Step 4: Let u_i be the i th column vector of K and initialize $\alpha_{i,j} = 1$.

Step 4.1: Repeat

Step 4.2: Compute β_j according to Eq. (29), i.e., $\beta_j = \sqrt{\sum_{i=1}^n \alpha_{i,j}^2 / \lambda}$.

Step 4.3: Compute α_i according to Eq. (27), i.e., $\alpha_i = (D_\beta^{-1} + \tilde{K})^{-1} \tilde{K}_i$.

Step 4.4: Until Convergence

Step 5: Rank the samples in X by following $\beta_j (j = 1, \dots, n)$ in a descending order and then return the top l samples as the selected most informative ones Z .

Output The l selected most informative samples can be labeled as the training samples.



Active Learning

$$\min_{w,L} \left\{ J(w) = \sum_{i=1}^m (w^T x_i - y_i)^2 + \mu \|w\|^2 + \gamma f^T L f \right\}$$

$$\min_{w,L} \{ J(w) = (X^T w - y)^2 + \mu \|w\|^2 + \gamma w^T X L X^T w \}$$

$$\min_{\mathbf{w}} \left\{ J(\mathbf{w}) = \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \mu \|\mathbf{w}\|^2 \right\}$$



$$\mathbf{C}_w = \left(\frac{\partial J(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right)^{-1} = (\mathbf{X}^\top \mathbf{X} + \mu \mathbf{I})^{-1}$$

$$\mathbf{C}_f = \mathbf{V} \mathbf{C}_w \mathbf{V}^\top$$



Diversifying Convex Transductive Experimental Design for Active Learning

12-IJCAI 2016

9-SIGIR 2008

$$\begin{aligned} \min_{\mathbf{A}, \mathbf{b}} \quad & \|\mathbf{X} - \mathbf{XA}\|_F^2 + \sum_{i=1}^n \frac{\sum_{j=1}^n a_{ij}^2}{b_i} + \gamma \|\mathbf{b}\|_1 + \alpha \mathbf{b}^T \mathbf{S} \mathbf{b} \\ \text{s.t.} \quad & b_i \geq 0, i = 1, 2, \dots, n \end{aligned}$$

$$\begin{aligned} \min_{\boldsymbol{\beta}, \boldsymbol{\alpha}_i \in \mathbb{R}^N} \quad & \sum_{i=1}^M \|\mathbf{x}_i - \mathbf{X}_C^T \boldsymbol{\alpha}_i\|^2 + \sum_{j=1}^N \frac{\alpha_{i,j}^2}{\beta_j} + \gamma \|\boldsymbol{\beta}\|_1 \\ \text{subject to} \quad & \mathbf{x}_i \in \mathbf{X}_{\mathcal{P}}, \beta_j \geq 0, j = 1, \dots, N, \end{aligned}$$

This constraint guarantees that highly similar samples would not have higher scores in sample selection at the same time.